# Monitoring Lustre I/O on the Franklin Cray XT4

**Andrew Uselton**

**NERSC**

**acuselton@lbl.gov**

# I/O on Franklin

- **Early Franklin file system performance was not where we wanted it**
  - **Interactive responsiveness at the command line**
  - **Bandwidth in parallel I/O**
- **Things have gotten better**
  - **LDAP changes improved responsiveness**
  - **We are better able to understand bandwidth issues**
- **We are looking for ways make things even better**
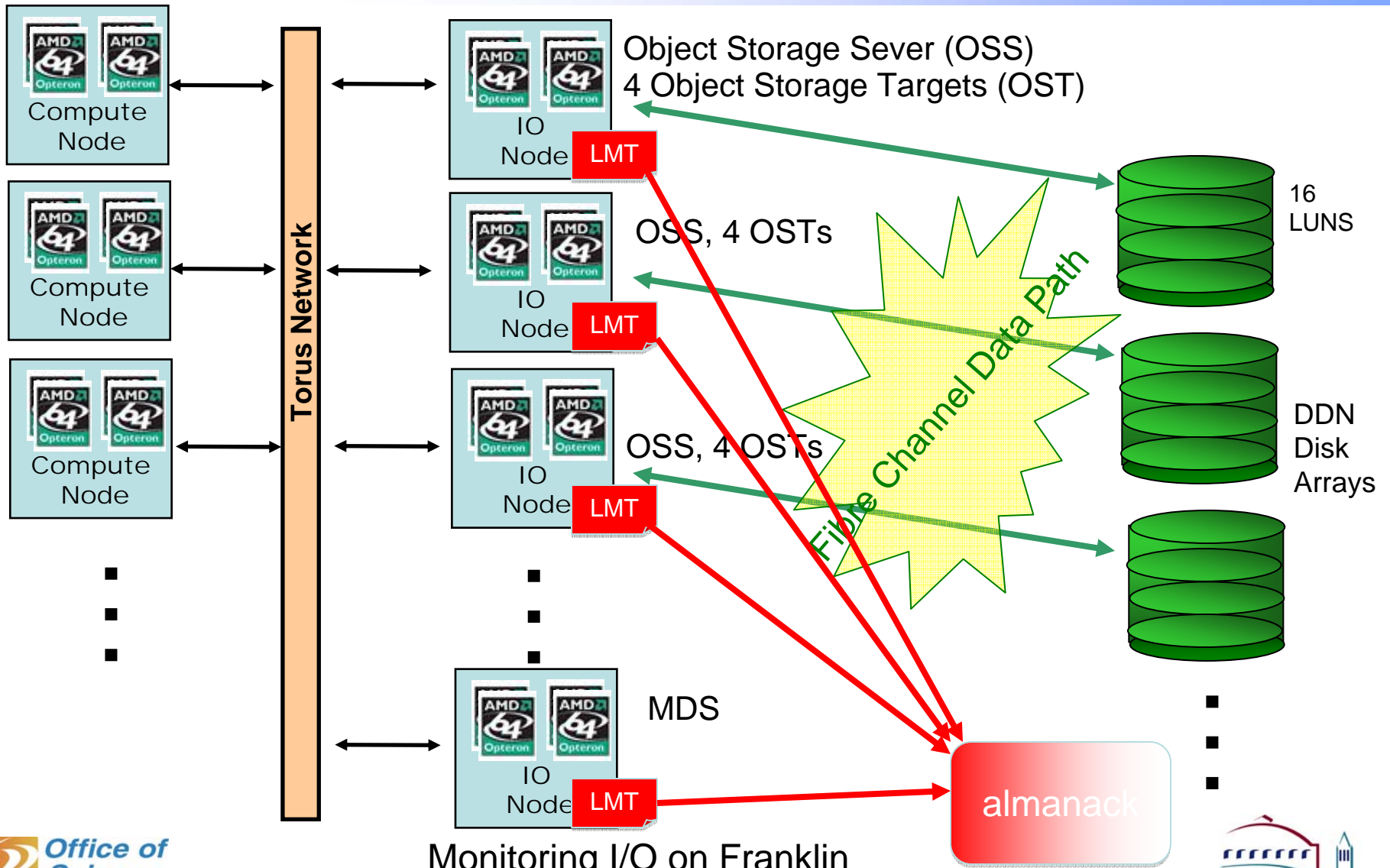  - **Monitoring**

# The Uses of Monitoring

- **After-the-fact incident investigation**
- **Checking the health of the system**
- **Noticing and exploring anomalous behavior**
- **Documenting use statistics**
- **System analysis and tuning**

# Overview

- **The Franklin I/O Pipeline**
- **The Lustre Monitoring Tool (LMT)**
- **Accessing LMT data**
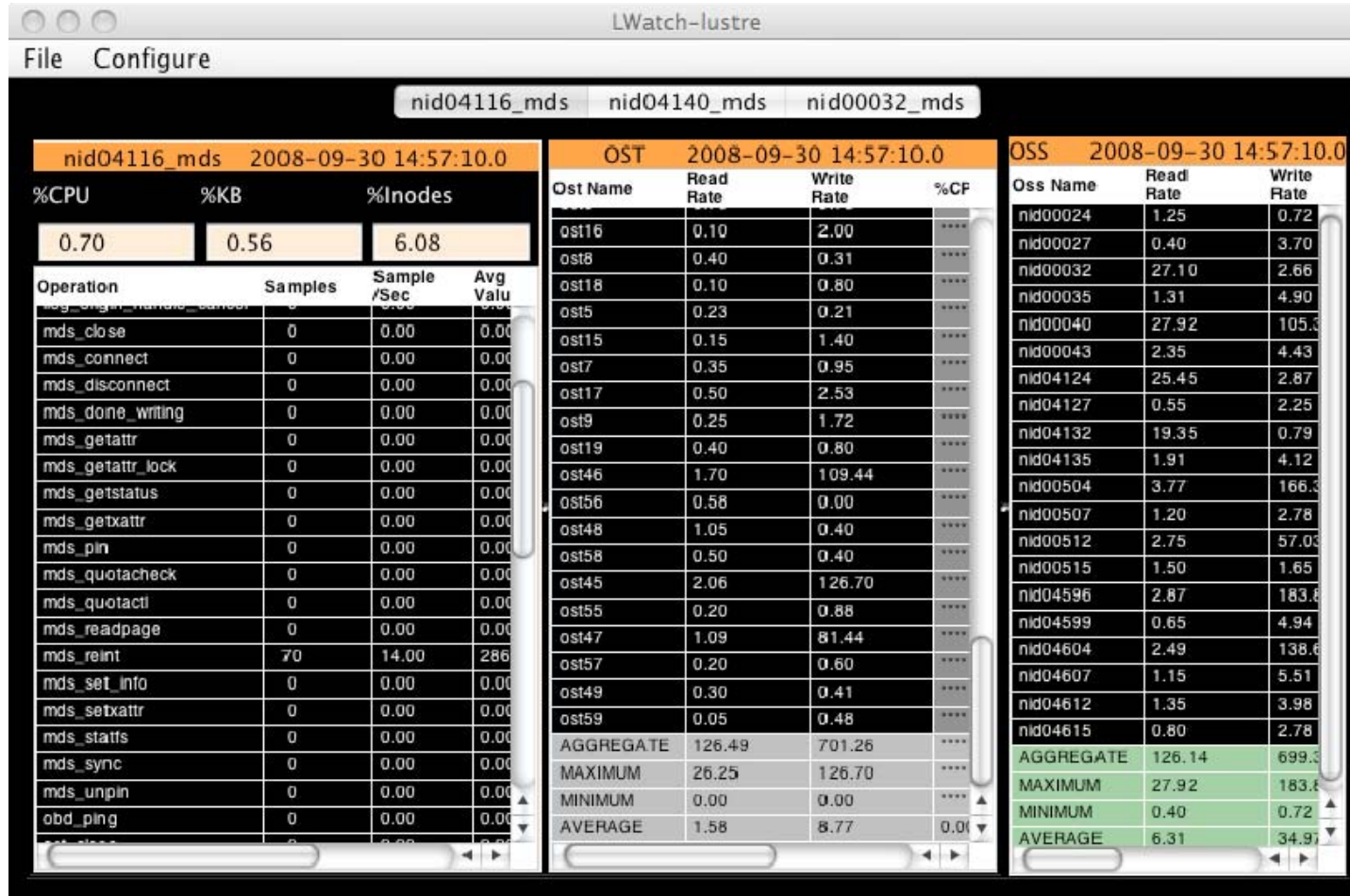- **Analyzing I/O performance**

# The Franklin I/O Pipeline

**Compute Node**

**Compute Node**

**Compute Node**

**Torus Network**

**IO Node** LMT

Object Storage Sever (OSS)
4 Object Storage Targets (OST)

**IO Node** LMT

OSS, 4 OSTs

**IO Node** LMT

OSS, 4 OSTs

**IO Node** LMT

MDS

Fibre Channel Data Path

16 LUNS

DDN Disk Arrays

almanack

Monitoring I/O on Franklin
5

# The Lustre Monitoring Tool

**The Lustre Monitoring Tool (LMT) is a set of plug-in modules for Cerebro, which is a monitoring package (resembling Ganglia). LMT and Cerebro were developed at LLNL.**

**LMT gathers data from /proc on Lustre servers and ships it to a database.**

**Additional tools query the data from there.**

Monitoring I/O on Franklin

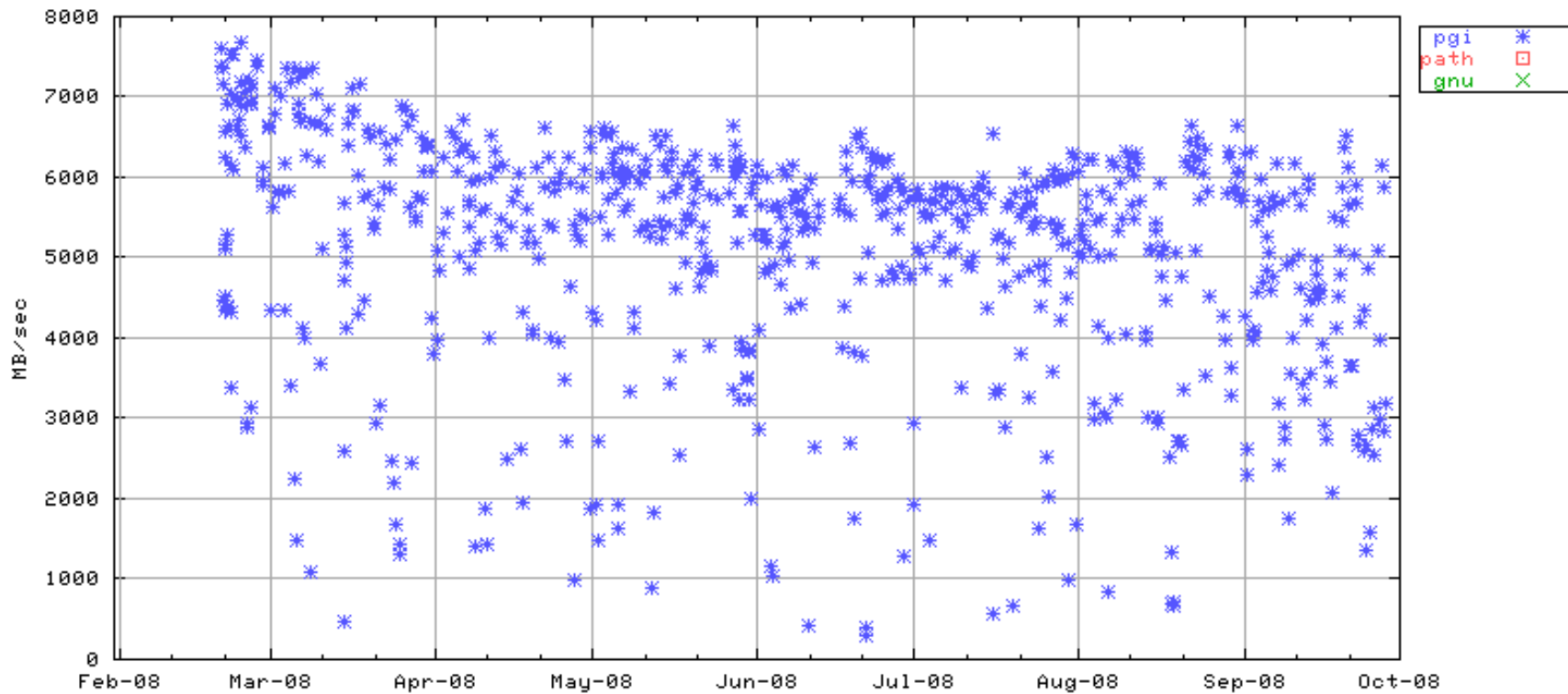# The LMT Dashboard

# Detail from the Dashboard

| | | |
|---|---|---|
| ost57 | 0.20 | 0.60 |
| ost49 | 0.30 | 0.41 |
| ost59 | 0.05 | 0.48 |
| AGGREGATE | 126.49 | 701.26 |
| MAXIMUM | 26.25 | 126.70 |
| MINIMUM | 0.00 | 0.00 |
| AVERAGE | 1.58 | 8.77 |

# LMT from the Command Line

bash-3.2$ /usr/share/lmt/bin/lmt.pl -s "2008-10-02 08:00:00" -e "2008-10-02 08:10:00"



Aggregate OST rates from 2008-10-02 08:00:00 to 2008-10-02 08:10:00

# Recent IOR Test Results (reads)



Monitoring I/O on Franklin

# A Particularly Slow IOR

# What Happened During My Job?

# Accessing LMT Data on the Web

# The View from LMT



Aggregate OST rates from 2008-09-25 04:03:00 to 2008-09-25 04:28:00

Monitoring I/O on Franklin

# A "Clean" IOR Run



Aggregate OST rates from 2008-09-12 00:12:00 to 2008-09-12 00:25:00

Monitoring I/O on Franklin

# 24 hours of I/O



Aggregate OST rates from 2008-09-29 08:00:00 to 2008-09-30 08:00:00

Monitoring I/O on Franklin

# Zooming in on the Reads



Aggregate OST rates from 2008-09-29 08:16:00 to 2008-09-29 08:46:00

Monitoring I/O on Franklin

17

# Useful Links

- **Franklin LMT web interface (may change)**

  http://www.nersc.gov/nusers/systems/franklin/lmt/

- **Franklin benchmarking page**

  https://www.nersc.gov/nusers/systems/franklin/monitor.php

- **LMT source code**

  http://sourceforge.net/projects/lmt/

- **Cerebro source code**

  http://sourceforge.net/projects/cerebro/